

This listing of claims will replace all prior versions, and listings, of claims in the application:

Claims 1-16 (canceled)

1 Claim 17 (currently amended): A system for building a
2 lexicon for use in capitalization correction for
3 unstructured excerpts, comprising:

4 a ripper adapted to assemble ~~assembling~~ a list of
5 word sets from unstructured content, ~~each word set at~~
6 least one of the word sets comprising a word and at
7 least two ~~one~~ variation on non-standard capitalization
8 variations for the word; and

9 an aggregator adapted to aggregate ~~aggregating~~
10 each word set, the aggregator including ~~comprising~~;

11 an analyzer adapted to identify ~~identifying~~
12 ~~at least one word set comprising significant~~
13 ~~statistics~~ non-standard capitalization variations
14 based on at least one criteria; and

15 a non-standard capitalization selector
16 adapted to select ~~selecting~~ at least ~~two such one~~
17 of the identified non-standard capitalization
18 variations within the identified word set ~~having~~
19 ~~a non-standard capitalization~~, and adding the at
20 ~~least two such~~ selected at least one of the
21 identified non-standard capitalization variations
22 to the lexicon, wherein the lexicon includes
23 records, each record including a word, wherein
24 the lexicon is indexed by the words included n
25 the records, and wherein at least one of the

26 records includes more than one non-standard
27 capitalization variation.

1 Claim 18 (currently amended): A system according to
2 Claim 17, further comprising:
3 a tokenizer adapted to tokenize ~~tokenizing~~ the
4 excerpt into the one or more words and one or more
5 punctuation marks.

1 Claim 19 (original): A system according to Claim 18,
2 wherein hyphenated words are split into a plurality of
3 the words.

Claim 20 (canceled)

1 Claim 21 (currently amended): A system according to
2 Claim ~~17~~ 20, wherein at least one of the non-standard
3 capitalization ~~comprises the at least one variation~~
4 ~~occurring~~ variations occurs in an excerpt having fewer
5 than half of individual letters provided in uppercase.

1 Claim 22 (currently amended): A system according to
2 Claim 17, further comprising:
3 a normalizer adapted to normalize ~~normalizing~~ a
4 plurality of the words extracted relative to a source
5 of the unstructured excerpt.

1 Claim 23 (currently amended): A system according to
2 Claim 17, wherein non-standard capitalization

3 variations that are identified based on one or more
4 criteria comprise the set comprising significant
5 statistics comprises only those non-standard
6 capitalization variations having at least four
7 occurrences ~~of at least one such variation within a~~
8 ~~word set.~~

1 Claim 24 (currently amended): A system according to
2 Claim 17, wherein at least one of the non-standard
3 capitalization ~~comprises the at least one variation~~
4 having variations has any individual letter other than
5 the first individual letter provided in uppercase.

Claim 25 (canceled)

1 Claim 26 (currently amended): A system according to
2 Claim 17, further comprising:
3 a validator adapted to apply ~~applying~~ implicit
4 rules for capitalization, and skipping each of the ~~at~~
5 ~~least two~~ non-standard capitalization variations
6 subject to at least one such implicit rule.

1 Claim 27 (currently amended): A system according to
2 Claim 26, wherein the implicit rules comprise skipping
3 each of the ~~at least two~~ non-standard capitalization
4 variations based on position within a sentence or
5 phrase.

1 Claim 28 (currently amended): A system according to
2 Claim 26, wherein the implicit rules comprise at least
3 one of (A) the non-standard capitalization variation
4 being a number, (B) the non-standard capitalization
5 variation having no vowels, and (C) the non-standard
6 capitalization variation constituting at least one of
7 an article, conjunction and preposition.

1 Claim 29 (currently amended): A system according to
2 Claim 26, wherein the implicit rules comprise
3 normalizing a number of occurrences for each of the
4 non-standard capitalization ~~at least two~~ variations
5 ~~using at least one of a normalizing function and~~
6 relative to a source of the non-standard
7 capitalization ~~each of the at least two~~ variations.

1 Claim 30 (currently amended): A system according to
2 Claim 26, wherein ~~the implicit rules comprise~~
3 ~~accommodating multiple forms of capitalization for~~
4 ~~each of the at least two variations by annotating each~~
5 ~~capitalization form~~ each of the word sets includes a
6 word and at least one non-standard capitalization
7 variation, each of the at least one non-standard
8 capitalization variation including with a frequency of
9 occurrence count ~~and skipping those of the each of the~~
10 ~~at least two variations occurring infrequently.~~

1 Claim 31 (original): A system according to Claim 17,
2 further comprising:
3 a hash table maintaining the lexicon.

1 Claim 32 (currently amended): A system according to
2 Claim 31, ~~further comprising:~~
3 ~~at least one record specifying at least one such~~
4 ~~word as a key into wherein the hash table is indexed~~
5 ~~by words , and associating at least one such variation~~
6 ~~within the word set as a preferred capitalization.~~

Claims 33-50 (canceled)

1 Claim 51 (new): A method comprising:
2 a) generating a plurality of word sets from a text
3 corpus, each of the words sets including
4 - a word identified from the text corpus,
5 - at least one capitalization variation, and
6 - a frequency of occurrence of each of the at
7 least one capitalization variation; and
8 b) generating a lexicon using the generated
9 plurality of word sets, wherein the lexicon
10 includes, for each of a plurality of words, at least
11 one capitalization variation identified using at
12 least one criteria, wherein at least one of the
13 words of the lexicon includes more than one
14 capitalization variation identified using the at
15 least one criteria.

1 Claim 52 (new): The method of claim 51 wherein a
2 capitalization variation is identified using the at least

3 one criteria only if it occurs at least four times in the
4 text corpus.

1 Claim 53 (new): The method of claim 51 further
2 comprising:

- 3 c) accepting a word having a capitalization
4 defining which, if any, of the characters of the
5 word are capitalized; and
- 6 d) performing a capitalization correction function
7 on the word using the generated lexicon.

1 Claim 54 (new): The method of claim 53 wherein the act
2 of performing a capitalization correction function
3 includes

- 4 - determining if the capitalization of the
5 word matches a capitalization variation in the
6 lexicon, and
- 7 - not changing the capitalization of the word
8 if it was determined to match a capitalization
9 variation in the lexicon.

1 Claim 55 (new): The method of claim 53 wherein the act
2 of performing a capitalization correction function
3 includes

- 4 - determining if the capitalization of the
5 word matches a capitalization variation in the
6 lexicon, which capitalization variation meets a
7 frequency criteria, and
- 8 - not changing the capitalization of the word
9 if it was determined to match a capitalization
10 variation in the lexicon.

1 Claim 56 (new): Apparatus comprising:
2 a) means for generating a plurality of word sets
3 from a text corpus, each of the words sets including
4 - a word identified from the text corpus,
5 - at least one capitalization variation, and
6 - a frequency of occurrence of each of the at
7 least one capitalization variation; and
8 b) means for generating a lexicon using the
9 generated plurality of word sets, wherein the
10 lexicon includes, for each of a plurality of words,
11 at least one capitalization variation identified
12 using at least one criteria, wherein at least one of
13 the words of the lexicon includes more than one
14 capitalization variation identified using the at
15 least one criteria.

1 Claim 57 (new): The apparatus of claim 56 wherein a
2 capitalization variation is identified using the at least
3 one criteria only if it occurs at least four times in the
4 text corpus.

1 Claim 58 (new): The apparatus of claim 56 further
2 comprising:
3 c) means for accepting a word having a
4 capitalization defining which, if any, of the
5 characters of the word are capitalized; and
6 d) means for performing a capitalization correction
7 function on the word using the generated lexicon.

1 Claim 59 (new): The apparatus of claim 58 wherein the
2 means for performing a capitalization correction function

3 - determine if the capitalization of the word
4 matches a capitalization variation in the
5 lexicon, and
6 - do not change the capitalization of the word
7 if it was determined to match a capitalization
8 variation in the lexicon.

1 Claim 60 (new): The apparatus of claim 58 wherein the
2 means for performing a capitalization correction function
3 - determine if the capitalization of the word
4 matches a capitalization variation in the
5 lexicon, which capitalization variation meets a
6 frequency criteria, and
7 - do not change the capitalization of the word
8 if it was determined to match a capitalization
9 variation in the lexicon.